

Supplementary webappendix

This webappendix formed part of the original submission and has been peer reviewed. We post it as supplied by the authors.

Supplement to: Milinovich GJ, Williams GM, Clements ACA, Hu W. Internet-based surveillance systems for monitoring emerging infectious diseases. *Lancet Infect Dis* 2013; published online Nov 28. [http://dx.doi.org/10.1016/S1473-3099\(13\)70244-5](http://dx.doi.org/10.1016/S1473-3099(13)70244-5).

Supplementary Table 1: Summary of publications describing internet-based surveillance systems for influenza or dengue (January 2008 – June 2013).

	Data source (resolution; range)^A	Region	Methodology	Results	Additional notes
Polgreen <i>et al.</i> (2008) ¹	Yahoo! search query logs (weekly; 2004-8)	United States	The authors correlated weekly aggregates of influenza-related search queries submitted to Yahoo! to national and regional influenza surveillance data. Linear models were applied with a 1-10 week lead time to predict influenza positive cultures and mortality rates attributable to pneumonia and influenza.	Correlations between estimates and culture results were reported to be $R^2=0.46$ (1 week lag) at a national level and $R^2=0.38$ (range =0.17-0.57; 1-3 weeks lag) regionally. Comparisons were also made with mortality data. At a national level, $R^2=0.42$ (5 weeks lag), while at a regional level $R^2=0.30$ (range =0.12-0.43; 4-6 weeks lag).	
Carneiro and Mylonakis (2009) ²	GI4S (weekly; 2004-9)	Worldwide and United States	The authors analysed GT search frequency for the term “bird flu”, worldwide and originating from the US between January 2004 and March 2009.	GT indicated an increase in search frequency for “bird flu” between 2005 and 2006 coinciding with the avian influenza outbreak; no cases recorded in the US.	The authors speculated that the spike in US search queries was driven by media reports.
Ginsberg <i>et al.</i> (2009) ³	Google (weekly; 2003-8)	United states	An automated method was used to identify Google search queries (2003-7) with the highest level of correlation with Center for Disease Control and Prevention (CDC) influenza-like illness (ILI) data. Models were produced using combinations of search terms and assessed against CDC ILI data excluded from the modeling process (2007-8).	The final model used 45 search queries and exhibited a high degree of correlation with ILI data; mean $r=0.90$ (range=0.80-0.96 for the 9 regions) for data used in the modeling process and $r=0.97$ (range =0.92-0.99) for hold-out data. Model estimates consistently preceded CDC surveillance reports by 1-2 weeks.	The model developed in this publication is utilised by GFT.
Hulth <i>et al.</i> (2009) ⁴	Vårdguiden (weekly; 2005-7)	Sweden	A model for estimating influenza incidence was created using frequency of queries submitted to Vårdguiden, a Swedish medical website. The model was validated against laboratory results and ILI reported by sentinel general practitioner (GP) clinics.	A model with four components was determined to be optimal for both sentinel ($R^2=0.89$; mean predictive error =0.08) and laboratory ($R^2=0.90$; mean predictive error =19.14) models.	This study provides proof-of-concept for monitoring Swedish search terms as a means of tracking infectious disease incidence.
Kelly and Grant (2009) ⁵	GFT (weekly; 2009)	Australia; Victoria	GFT data for Australia was compared to Victorian GP sentinel surveillance and Melbourne Medical Deputising Service data for ILI-related consultations during the 2009 H1N1 influenza pandemic.	GFT was reported to have “remarkable correlation” with both GP sentinel surveillance and Melbourne Medical Deputising Service data.	Statistical analyses are not reported.
Pattie <i>et al.</i> (2009) ⁶	GFT (weekly; 2007-8)	United States	Weekly influenza data from the Department of Defense’s Electronic Surveillance System for the Early Notification of Community-Based Epidemics (ESSENCE) were compared with GFT and CDC Outpatient Influenza-like Illness Network (ILINet) results.	Correlations between ESSENCE and CDC ILINet and GFT were $r=0.92$ and 0.88, respectively.	
Pelat <i>et al.</i> (2009) ⁷	GI4S (weekly; 2004-9)	France	Search query frequency for French terms related to ILI, gastroenteritis and chickenpox were compared with French Sentinel Network surveillance data.	The query “grippe –aviaire –vaccine” (influenza –avian –vaccine) produced the highest correlation with surveillance data for influenza ($r=0.82$; 0 week lag).	This study provides proof-of-concept for monitoring French search terms as a means of tracking infectious disease incidence.
Wilson <i>et al.</i> (2009) ⁸	GFT (weekly; 2009)	New Zealand	GFT data were compared graphically to ILI data from two independent national surveillance systems and with ILI data from a national toll-free telephone triage and health service.	GFT exhibited good visual correlation with the data from the two national surveillance systems (peaked within one week). Peak ILI-related calls to Healthline occurred three weeks prior to GFT.	
Chew and Eysenbach (2010) ⁹	Twitter (Weekly; 2009)	United States	Twitter posts containing the terms “swine flu”, “swineflu” and/or “H1N1” were collected and automatically binned into groups describing the nature of the content. The proportion and absolute numbers of posts binned as “personal experiences” or “concern” were compared to US H1N1 incidence rates.	Correlations between H1N1 incidence rates and posts binned as “personal experiences” were reported to be $r=0.77$ for absolute number of posts and $r=0.67$ for the percentage of posts assigned to this category. Correlations for posts assigned to the “concern” category were $r=0.66$ for absolute number and $r=0.37$ for the percentage of posts.	
Corley <i>et al.</i>	Spinn3r	United States	The authors used Spinn3r used to monitor blogs for terms	Level of correlation between flu-content posts and	Note: Corley <i>et al.</i> (2010) ¹⁰ and

Supplementary webappendix

This webappendix formed part of the original submission and has been peer reviewed. We post it as supplied by the authors.

Supplement to: Milinovich GJ, Williams GM, Clements ACA, Hu W. Internet-based surveillance systems for monitoring emerging infectious diseases. *Lancet Infect Dis* 2013; published online Nov 28. [http://dx.doi.org/10.1016/S1473-3099\(13\)70244-5](http://dx.doi.org/10.1016/S1473-3099(13)70244-5).

(2010) ¹⁰	(weekly; 2008-9)		“influenza” and “flu”. This was correlated with ILINet data (October 2008 and January 2009).	ILINet data was reported to be $r=0.626$.	Corley <i>et al.</i> (2010) ¹¹ report the same dataset. Length of analyses period differs.
Corley <i>et al.</i> (2010) ¹¹	Spinn3r (weekly; 2008-9)	United States	he authors used Spinn3r used to monitor blogs for terms “influenza” and “flu”. This was correlated with ILINet data (October 2008 and March 2009).	Level of correlation between flu-content posts and ILINet data was reported to be $r=0.545$.	Note: Corley <i>et al.</i> (2010) ¹⁰ and Corley <i>et al.</i> (2010) ¹¹ report the same dataset. Length of analyses period differs.
Valdivia <i>et al.</i> (2010) ¹²	GFT (weekly; 2009)	Belgium, Bulgaria, France, Germany, Hungary, Netherlands, Norway, Poland, Russian Federation, Spain, Sweden, Switzerland, and Ukraine	Correlations between GFT and sentinel network data over the course of the 2009 influenza A (H1N1) pandemic were assessed. Data were analysed over three periods: the periods before and after 31 August, 2009 (pre- and during the epidemic) and for the entire period (23 March, 2009 to 28 March, 2010).	Overall correlation coefficients were reported to be $r=0.72-0.94$. Pre-epidemic correlations were $r=0.39-0.87$ and correlations during the epidemic period were $r=0.53-0.97$. GFT peak incidence occurred 0-2 weeks prior to sentinel estimates, except for Sweden (GFT preceded sentinel network estimates by 11 weeks).	
Valdivia and Monge-Corella (2010) ¹³	GI4S (weekly; 2004-9)	Spain	Search query frequency for Spanish terms related to ILI and chickenpox were compared with data from the Spanish National Epidemiology Center.	The query “gripe – aviar – vacuna” (influenza –avian – vaccine) produced the highest correlation with surveillance data for influenza ($r=0.81$; 2 week lag).	This study provides proof-of-concept for monitoring Spanish search terms as a means of tracking infectious disease incidence.
Althouse <i>et al.</i> (2011) ¹⁴	GI4S (weekly or monthly; 2004-11)	Singapore and Bangkok	Search frequencies for dengue-related terms were downloaded from GI4S and categorised as: nomenclature, signs/symptoms and treatment. Three models for predicting incidence were applied to the data (step-down linear regression, generalized boosted regression, and negative binomial regression) and performance assessed against surveillance data from the Singapore Ministry of Health (weekly) and the Thai Bureau of Epidemiology (monthly). Additionally, Support Vector Machine (SVM) and logistic regression models were applied to the data to predict periods of high incidence (thresholds: 50 th , 75 th and 90 th percentile of case numbers).	Step-down linear regression provided the closest fit to the surveillance data (Singapore $R^2=0.948$, $n=16$ search terms; Bangkok ($R^2=0.943$, $n=8$ search terms). SVM was superior to logistic regression in determining periods of high incidence. Using the 75 th percentile cut off, the area under the receiver operating characteristic curve was 0.91 and 0.96 for Singapore and Bangkok, respectively.	
Boyle <i>et al.</i> (2011) ¹⁵	GFT (weekly; 2006-9)	Australia; Queensland	The authors compared historical data for presentation and admissions to 27 Queensland public hospitals, for ILI, during the influenza season (June-September) with corresponding GFT data.	Correlation coefficients for the GFT data with ILI data were $r=0.35$; 0.88; 0.91; 0.76; respectively for the 2006-2009 influenza seasons.	
Chan <i>et al.</i> (2011) ¹⁶	Google (weekly or monthly; 2003-10)	Bolivia, Brazil, India, Indonesia and Singapore	A model for dengue was created using a methodology similar to that described above for GFT. ³ Weekly (Bolivia and Singapore) or monthly (Brazil, India and Indonesia) Google search query aggregates were used to create individual models and performance was assessed against national surveillance data. The models were engineered to remove spurious spikes and the validated using holdout datasets.	Models utilised up to ten search terms. Estimates produced were described to match observed curves “quite well” for all countries. Pearson’s correlation coefficients for the overall and holdout data sets, respectively, were: Bolivia ($r=0.94$, 0.83), Brazil ($r=0.92$, 0.99), India ($r=0.87$, 0.94), Indonesia ($r=0.90$, 0.94) and Singapore ($r=0.82$, 0.94).	The model developed in this publication is utilised by Google Dengue Trends.
Collier <i>et al.</i> (2011) ¹⁷	Twitter (weekly; 2009-10)	United States	Supervised learning was used to categorise 97 million Twitter messages into five influenza-related categories. Weekly volumes were compared against US CDC data for positive Influenza tests.	Moderately strong correlations ($r=0.58-0.67$) were observed between positive message frequency in each class and CDC results for Influenza A (H1N1).	The authors conceded that further progress needs to be made to achieve the high degree of correlation achieved by GFT.
Cook <i>et al.</i> (2011) ¹⁸	GFT (weekly; 2009)	United States	A new GFT model was developed using a larger pool of candidate queries and more relaxed parameters than the original model. ³ Both models were compared with ILINet estimates; data were analysed over four time periods covering the 2009 H1N1 outbreak (pre-H1N1, Summer	Both models exhibited a high level of correlation pre-H1N1, during Winter and over the entire period analysed ($r>0.90$). Correlations for the Summer H1N1 period were higher for the new model ($r = 0.95$ vs. 0.29).	Internet search behaviour changed during the H1N1 pandemic. The new GFT model, trained using data from the summer months of the pandemic, was able to accommodate

			H1N1, Winter H1N1, and H1N1 overall).		changes in search behaviour.
Hulth and Rydevik (2011) ¹⁹	GI4S, GFT & Vårdguiden (weekly; 2009)	Sweden	Search frequencies for the term “influenza” submitted to Google (using GI4S) and the Vårdguiden web site were compiled and visually compared to estimates for Sweden produced by GFT, the Vårdguiden model and Swedish sentinel reports.	The authors reported the Vårdguiden model to produce the best representation of the sentinel reports.	This publication presents a system (Generating Epidemiological Trends from Web Logs, Like; GETWELL) which integrates automated, internet-based surveillance for influenza into an existing traditional surveillance framework.
Hulth and Rydevik (2011) ²⁰	GFT & Vårdguiden (weekly; 2009)	Sweden	The Vårdguiden model ⁴ and GFT (Sweden) performance was analysed over the 2009 H1N1 outbreak. Models were compared with “incomplete sentinel” (weekly reports) and “complete sentinel” (final reports distributed five weeks later, containing late reports) data.	Correlation between the Vårdguiden model and complete sentinel data ($r=0.90$, $R^2=0.75$) were higher than for the GFT model ($r=0.87$, R^2 not reported). A similar trend was reported for incomplete sentinel data: Vårdguiden model ($r=0.88$, $R^2=0.68$), GFT model ($r=0.85$, R^2 not reported).	
Malik <i>et al.</i> (2011) ²¹	GFT (weekly; 2009)	Canada; Manitoba	Performance of GFT (Manitoba) was evaluated during the 2009 H1N1 pandemic using laboratory confirmed H1N1 cases as a reference. Correlations between emergency department (ED) indicators and virological data were also determined.	GFT exhibited highest correlation with virological data when a 2-week lag was incorporated ($R^2=0.69$). GFT also exhibited a high level of correlation with ED indicators analysed ($R^2=0.86$ and $R^2=0.85$).	
Ortiz <i>et al.</i> (2011) ²²	GFT (weekly; 2003-8)	United States	Surveillance data from the US Influenza Virologic Surveillance System was used as a reference to analyse performance of GFT and the ILINet. A secondary analysis was performed to determine the influence of high-leverage outlier points.	GFT correlated highly with CDC ILI surveillance data ($r=0.94$) and CDC Virus Surveillance ($r=0.72$). Removal of outliers increased correlation of GFT with CDC Virus Surveillance ($r=0.82$), but not CDC ILINet data ($r=0.72$). GFT and CDC Virus Surveillance correlations differed by season ($r=0.67$ - 0.94 , mean $=0.79$) and region ($r=0.64$ - 0.80 , mean $=0.70$).	GFT and CDC ILI surveillance data shared 88% of variance, whereas, GFT and CDC Virus Surveillance shared 51%.
Dugas <i>et al.</i> (2012) ²³	GFT (weekly; 2009-10)	United States; Baltimore, Maryland	City GFT data for Baltimore were correlated with separate adult and paediatric data for ILI, laboratory confirmed influenza cases and ED crowding indices.	GFT exhibited a high degree of correlation with ED presentation for ILI (adult $r=0.89$; paediatric $r=0.65$) and laboratory-confirmed influenza cases (adult $r=0.88$; paediatric $r=0.72$). GFT and crowding metrics exhibited good correlation for paediatric ED visit volumes ($r=0.65$), and number of patients who left without consulting a clinician ($r=0.64$). Peak correlations (if present) were observed for all data analysed with 0-1 week lag.	
Dukic <i>et al.</i> (2012) ²⁴	GFT (Weekly; 2003-9)	United States	The authors used GFT data to track influenza dynamics using a season-specific susceptible-exposed-infected-recovered (SEIR) model within the state-space framework. The study also compared posterior distributions between Markov chain Monte Carlo (MCMC) algorithm and sequential learning algorithm. National and statewide data from nine states were used.	This publication demonstrated the ability of GFT data, applied to a state-space SEIR model, to provide near real-time disease tracking. The sequential learning algorithm could improve the computational speed in comparison with MCMC algorithm.	
Lamos and Cristianini (2012) ²⁵	Twitter (weekly; 2009-10)	United Kingdom; regional	The authors harvested a sample of Twitter posts tagged with locations identifiable to highly populated urban centres in the UK. Influenza incidence was estimated by monitoring frequency of Twitter posts containing key words selected using the methods of term frequencies and bootstrapped least absolute shrinkage and selection operator. Estimates were compared to official influenza incidence rates provided by the Health Protection Agency.	Linear correlation coefficients of up to $r=0.933$ were reported. The authors conclude that supervised learning framework presents a suitable method for selecting features for use in digital surveillance systems.	
Patwardhan	GFT (weekly; 2009-10)	United States	GFT results were correlated with sales of four drugs	Correlation between GFT and pharmaceuticals sales	

and Bilkovski (2012) ²⁶	2007-11)		commonly prescribed for influenza. Pharmaceutical sales were also correlated with 2007 ILINet data.	was $r=0.92$. Correlation between pharmaceutical sales and 2007 ILINet data was $r=0.97$.	
Pervaiz <i>et al.</i> (2012) ²⁷	GFT (weekly; 2003-11)	United States	The authors explored the potential of using GFT as a basis for building fully automated early warning systems. The accuracy and practicality of three types of algorithms (normal distribution, Poisson distribution and negative binomial distribution) were assessed.	Poisson and negative binomial regression models were demonstrated to perform better, on average, than those based on normal distribution. The authors concluded that the addition of a level of computational intelligence to GFT data provided a reliable means of early outbreak detection.	
Scarpino <i>et al.</i> (2012) ²⁸	GFT (weekly; 2001-8)	United States; Texas	The authors used a submodular optimisation algorithm to optimise sentinel provider selection for the Texas ILINet. The effect of incorporating GFT as a virtual provider into the network was investigated.	GFT alone matched the performance of an optimised network of 44 providers and outperformed the 2008 Texas ILINet (82 providers). Inclusion of GFT into an optimised network of 82 providers increased performance by 12.5%.	This publication demonstrates the potential for using GFT data to augment current surveillance systems.
Shaman and Karspeck (2012) ²⁹	GFT (Weekly; 2003-11)	United States; New York City, New York	The authors used the ensemble adjustment Kalman filter to assimilate GFT data and then applied a humidity-driven susceptible-infectious-recovered-susceptible model.	The authors reported that the approach used was able to provide indication of peak influenza incidence up to seven weeks in advance of its occurrence.	
Culotta (2013) ³⁰	Twitter (Weekly; 2009-10)	United States	Frequency of posts containing influenza-related keywords were analysed in a corpus of 570 million Twitter messages. Changes in frequency of posts containing these words (or combinations) were correlated with US CDC data for influenza. A document classification method was used to filter spurious messages and correlations assessed in the absence of these messages.	Linear models fitted to influenza related tweet frequency were demonstrated to exhibit strong correlations with CDC ILI hold out data (up to $r=0.92$ for a single term and $r=0.97$ for a combination of terms). The inclusion of a document classification method was reported to limit the effect of spurious posts.	
Dugas <i>et al.</i> (2013) ³¹	GFT (Weekly; 2004-11)	United States; a single urban, tertiary care emergency department (location undisclosed)	The authors developed an influenza forecast model that used data available in real-time at the city or medical center level. Forecast models were produced using a Negative Binomial Generalized Autoregressive Moving Average (GARMA) model. Secondary variables (GFT, meteorological data and temporal variables) were added as external variables and predictive performance assessed.	GFT data was demonstrated to outperform meteorological data and temporal variables in generalized linear models. The inclusion of GFT into the selected GARMA base model significantly improved performance (from 81% to 83%; $p=0.0005$). Inclusion of other variables to this model did not further increase performance.	This paper demonstrates proof-of-concept for use of GFT in forecasting models.
Kang <i>et al.</i> (2013) ³²	GT (Weekly; 2008-11)	China; Guangdong province	The authors compared GT indices for the terms “flu”, “common cold”, “fever”, “cough”, “sore throat”, “influenza A” and “H1N1” with Guangdong province ILI surveillance and laboratory confirmed influenza (virological) data (from the Guangdong CDC).	The highest level of correlation reported was between “fever” and ILI ($r=0.73$; no lag). “Influenza A” had the highest level of correlation with virological surveillance data ($r=0.66$; one week lag). Lag times of 0 weeks produced the highest levels of correlations for most search terms.	Correlation between ILI and virological surveillance data for Guangdong province were reported to be $r=0.56$ (95% CI: 0.43, 0.66).
Vandendijck <i>et al.</i> (2013) ³³	GFT (Weekly; 2003-11)	Belgium; Flanders	The authors compared ILI estimates in Flanders produced by the Great Influenza Survey (an online weekly influenza survey), GFT and the Belgian Sentinel Network.	Correlations between GFT and the Belgian Sentinel Network were not reported in this study. The Great Influenza Survey (random walk model), however, exhibited good correlation with both GFT ($r=0.62-0.94$) and Belgian Sentinel Network ($r=0.76-0.94$).	
Yuan <i>et al.</i> (2013) ³⁴	Baidu (Monthly; 2009-12)	China	The authors used search frequency data collected from Baidu (the most commonly used search engine in China) to estimate influenza activity in China. Estimates were compared to official influenza counts from the Ministry of Health (MOH). Additionally, the authors created a hybrid model (that uses both Baidu search index and MOH case report data) to produce estimates of influenza incidence.	The search index composite created by the authors for estimating influenza incidence contained 8 terms and exhibited a correlation with MOH data of $r=0.96$ (lag 0). The predictive model was reported to have an R^2 value of 0.95. Mean absolute percent error of this model on the eight month out-of-sample period analysed was 10.6% (1.8%-22.2%).	This is the first reported use of Baidu data for influenza surveillance.
Zhou <i>et al.</i> (2013) ³⁵	GT (Daily; 2006-11)	United States	Search query frequencies were used to predict epidemic alert levels for influenza by a continuous density Hidden Markov model (HMM). Queries used in this publication were	HMM-based methods were reported to predict influenza alert levels with 97.7% accuracy. Bayes and regression-based methods did so with 85.3% and 81.3% accuracy	The authors were able to produce estimates of influenza alert levels in “real-time” as GT search volumes

			selected using a greedy forward selection method. The accuracy of these estimates was compared to alert levels calculated from US CDC data.	respectively.	are updated daily.
--	--	--	---	---------------	--------------------

^A Google Flu Trends <http://www.google.org/flutrends/>; Google Insights for Search (merged with Google Trends in 2012); Google Trends <http://www.google.com/trends/>; Yahoo! <http://search.yahoo.com/>; Vårdguiden <http://www.vardguiden.se/>; Baidu <http://index.baidu.com/>

Acronyms

CDC	Center for Disease Control and Prevention
ED	Emergency Department
ESSENCE	Early Notification of Community-Based Epidemics
GARMA	Generalised Autoregressive Moving Average
GET WELL	Generating Epidemiological Trends from WEb Logs, Like
GFT	Google Flu Trends
GI4S	Google Insights for Search
GP	General Practitioner
GT	Google Trends
HMM	Hidden Markov Model
ILI	Influenza-like Illness
ILINet	Outpatient Influenza-like Illness Network
MCMC	Markov Chain Monte Carlo
MOH	Ministry of Health
SEIR	Susceptible-Exposed-Infected-Recovered
SVM	Support Vector Machine

References

1. Polgreen PM, Chen Y, Pennock DM, Nelson FD. Using internet searches for influenza surveillance. *Clinical infectious diseases : an official publication of the Infectious Diseases Society of America* 2008; **47**: 1443-8.
2. Carneiro HA, Mylonakis E. Google trends: a web-based tool for real-time surveillance of disease outbreaks. *Clinical infectious diseases : an official publication of the Infectious Diseases Society of America* 2009; **49**: 1557-64.
3. Ginsberg J, Mohebbi MH, Brammer L, Smolinski MS, Brilliant L. Detecting influenza epidemics using search engine query data. *Nature* 2009; **457**: 1012-4.
4. Hulth A, Rydevik G, Linde A. Web queries as a source for syndromic surveillance. *PloS one* 2009; **4**: e4378.
5. Kelly H, Grant K. Interim analysis of pandemic influenza (H1N1) 2009 in Australia: surveillance trends, age of infection and effectiveness of seasonal vaccination. *Euro surveillance : bulletin europeen sur les maladies transmissibles = European communicable disease bulletin* 2009; **14**.
6. Pattie DC, Cox KL, Burkom HS, Lombardo JS, Gaydos JC. A public health role for Internet search engine query data? *Military medicine* 2009; **174**: xi-xii.
7. Pelat C, Turbelin C, Bar-Hen A, Flahault A, Valleron A. More diseases tracked by using Google Trends. *Emerging infectious diseases* 2009; **15**: 1327-8.
8. Wilson N, Mason K, Tobias M, Peacey M, Huang QS, Baker M. Interpreting Google flu trends data for pandemic H1N1 influenza: the New Zealand experience. *Euro surveillance : bulletin europeen sur les maladies transmissibles = European communicable disease bulletin* 2009; **14**.
9. Chew C, Eysenbach G. Pandemics in the age of Twitter: content analysis of Tweets during the 2009 H1N1 outbreak. *PloS one* 2010; **5**: e14118.
10. Corley CD, Cook DJ, Mikler AR, Singh KP. Using Web and social media for influenza surveillance. *Advances in experimental medicine and biology* 2010; **680**: 559-64.

11. Corley CD, Cook DJ, Mikler AR, Singh KP. Text and structural data mining of influenza mentions in Web and social media. *International journal of environmental research and public health* 2010; **7**: 596-615.
12. Valdivia A, Lopez-Alcalde J, Vicente M, Pichiule M, Ruiz M, Ordobas M. Monitoring influenza activity in Europe with Google Flu Trends: comparison with the findings of sentinel physician networks - results for 2009-10. *Euro surveillance : bulletin europeen sur les maladies transmissibles = European communicable disease bulletin* 2010; **15**.
13. Valdivia A, Monge-Corella S. Diseases tracked by using Google trends, Spain. *Emerging infectious diseases* 2010; **16**: 168.
14. Althouse BM, Ng YY, Cummings DA. Prediction of dengue incidence using search query surveillance. *PLoS neglected tropical diseases* 2011; **5**: e1258.
15. Boyle JR, Sparks RS, Keijzers GB, Crilly JL, Lind JF, Ryan LM. Prediction and surveillance of influenza epidemics. *The Medical journal of Australia* 2011; **194**: S28-33.
16. Chan EH, Sahai V, Conrad C, Brownstein JS. Using web search query data to monitor dengue epidemics: a new model for neglected tropical disease surveillance. *PLoS neglected tropical diseases* 2011; **5**: e1206.
17. Collier N, Son NT, Nguyen NM. OMG U got flu? Analysis of shared health messages for bio-surveillance. *Journal of biomedical semantics* 2011; **2 Suppl 5**: S9.
18. Cook S, Conrad C, Fowlkes AL, Mohebbi MH. Assessing Google flu trends performance in the United States during the 2009 influenza virus A (H1N1) pandemic. *PloS one* 2011; **6**: e23610.
19. Hulth A, Rydevik G. GET WELL: an automated surveillance system for gaining new epidemiological knowledge. *BMC public health* 2011; **11**: 252.
20. Hulth A, Rydevik G. Web query-based surveillance in Sweden during the influenza A(H1N1)2009 pandemic, April 2009 to February 2010. *Euro surveillance : bulletin europeen sur les maladies transmissibles = European communicable disease bulletin* 2011; **16**.
21. Malik MT, Gumel A, Thompson LH, Strome T, Mahmud SM. "Google flu trends" and emergency department triage data predicted the 2009 pandemic H1N1 waves in Manitoba. *Canadian journal of public health Revue canadienne de sante publique* 2011; **102**: 294-7.
22. Ortiz JR, Zhou H, Shay DK, Neuzil KM, Fowlkes AL, Goss CH. Monitoring influenza activity in the United States: a comparison of traditional surveillance systems with Google Flu Trends. *PloS one* 2011; **6**: e18687.
23. Dugas AF, Hsieh YH, Levin SR, et al. Google Flu Trends: correlation with emergency department influenza rates and crowding metrics. *Clinical infectious diseases : an official publication of the Infectious Diseases Society of America* 2012; **54**: 463-9.
24. Dukic V, Lopes HF, Polson NG. Tracking Epidemics With Google Flu Trends Data and a State-Space SEIR Model. *Journal of the American Statistical Association* 2012; **107**: 1410-26.
25. Lamos V, Cristianini N. Nowcasting Events from the Social Web with Statistical Learning. *Acm Transactions on Intelligent Systems and Technology* 2012; **3**.
26. Patwardhan A, Bilkovski R. Comparison: Flu prescription sales data from a retail pharmacy in the US with Google Flu trends and US ILINet (CDC) data as flu activity indicator. *PloS one* 2012; **7**: e43611.
27. Pervaiz F, Pervaiz M, Abdur Rehman N, Saif U. FluBreaks: Early Epidemic Detection from Google Flu Trends. *Journal of medical Internet research* 2012; **14**: e125.
28. Scarpino SV, Dimitrov NB, Meyers LA. Optimizing provider recruitment for influenza surveillance networks. *PLoS computational biology* 2012; **8**: e1002472.
29. Shaman J, Karspeck A. Forecasting seasonal outbreaks of influenza. *Proc Natl Acad Sci U S A* 2012; **109**: 20425-30.
30. Culotta A. Lightweight methods to estimate influenza rates and alcohol sales volume from Twitter messages. *Language Resources and Evaluation* 2013; **47**: 217-38.
31. Dugas AF, Jalalpour M, Gel Y, et al. Influenza forecasting with Google Flu Trends. *PloS one* 2013; **8**: e56176.
32. Kang M, Zhong H, He J, Rutherford S, Yang F. Using Google Trends for influenza surveillance in South China. *PloS one* 2013; **8**: e55205.
33. Vandendijck Y, Faes C, Hens N. Eight years of the great influenza survey to monitor influenza-like illness in Flanders. *PloS one* 2013; **8**: e64156.
34. Yuan Q, Nsoesie EO, Lv B, Peng G, Chunara R, Brownstein JS. Monitoring influenza epidemics in China with search query from Baidu. *PloS one* 2013; **8**: e64323.
35. Zhou XC, Li Q, Zhu ZL, Zhao H, Tang H, Feng YJ. Monitoring Epidemic Alert Levels by Analyzing Internet Search Volume. *Ieee Transactions on Biomedical Engineering* 2013; **60**: 446-52.